

Linnea: Automatic Generation of Efficient Linear Algebra Programs

Henrik Barthels

Introduction

- How to compute the following expressions?

$$\mathbf{b} := (\mathbf{X}^T \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}^{-1} \mathbf{y}$$

$$\mathbf{x} := \mathbf{W}(\mathbf{A}^T (\mathbf{A} \mathbf{W} \mathbf{A}^T)^{-1} \mathbf{b} - \mathbf{c})$$

$$\mathbf{x} := (\mathbf{A}^{-T} \mathbf{B}^T \mathbf{B} \mathbf{A}^{-1} + \mathbf{R}^T [\Lambda(\mathbf{R} \mathbf{z})] \mathbf{R})^{-1} \mathbf{A}^{-T} \mathbf{B}^T \mathbf{B} \mathbf{A}^{-1} \mathbf{y}$$

$$\mathbf{X}_{k+1} := \mathbf{S}(\mathbf{S}^T \mathbf{A} \mathbf{S})^{-1} \mathbf{S}^T + (\mathbf{I}_n - \mathbf{S}(\mathbf{S}^T \mathbf{A} \mathbf{S})^{-1} \mathbf{S}^T \mathbf{A}) \mathbf{X}_k (\mathbf{I}_n - \mathbf{A} \mathbf{S}(\mathbf{S}^T \mathbf{A} \mathbf{S})^{-1} \mathbf{S}^T)$$

- High-level languages are easy to use, but performance is usually suboptimal.
- BLAS and LAPACK can be fast, but require a lot of expertise.
- Goal: Simplicity **and** performance.
- Dense, mid- to large-scale linear algebra.

Introduction

How to compute...

$$y' := H^\dagger y + (I_n - H^\dagger H)x \quad [\text{TG17}]$$

...with these operations?

$$x := Ab \quad 2n^2$$

$$X := AB \quad 2n^3$$

$$x := a \pm b \quad n$$

$$X := A \pm B \quad n^2$$

Introduction

How to compute...

$$y' := H^\dagger y + (I_n - H^\dagger H)x \quad [\text{TG17}]$$

...with these operations?

$$x := Ab \quad 2n^2$$

$$X := AB \quad 2n^3$$

$$x := a \pm b \quad n$$

$$X := A \pm B \quad n^2$$

$$M_1 := H^\dagger H$$

$$M_2 := I_n - M_1$$

$$m_3 := M_2 x$$

$$m_4 := H^\dagger y$$

$$y' := m_3 + m_4$$

$$\Rightarrow 2n^3 + 5n^2 + n \text{ FLOPs}$$

Introduction

How to compute...

$$y' := H^\dagger y + (I_n - H^\dagger H)x \quad [\text{TG17}]$$

$$\Leftrightarrow y' := H^\dagger y + x - H^\dagger Hx$$

$$\Leftrightarrow y' := H^\dagger (y - Hx) + x$$

...with these operations?

$$x := Ab \quad 2n^2$$

$$X := AB \quad 2n^3$$

$$x := a \pm b \quad n$$

$$X := A \pm B \quad n^2$$

$$M_1 := H^\dagger H$$

$$M_2 := I_n - M_1$$

$$m_3 := M_2 x$$

$$m_4 := H^\dagger y$$

$$y' := m_3 + m_4$$

$$\Rightarrow 2n^3 + 5n^2 + n \text{ FLOPs}$$

Introduction

How to compute...

$$y' := H^\dagger y + (I_n - H^\dagger H)x \quad [\text{TG17}]$$

$$\Leftrightarrow y' := H^\dagger y + x - H^\dagger Hx$$

$$\Leftrightarrow y' := H^\dagger (y - Hx) + x$$

...with these operations?

$$x := Ab \quad 2n^2$$

$$X := AB \quad 2n^3$$

$$x := a \pm b \quad n$$

$$X := A \pm B \quad n^2$$

$$M_1 := H^\dagger H$$

$$M_2 := I_n - M_1$$

$$m_3 := M_2 x$$

$$m_4 := H^\dagger y$$

$$y' := m_3 + m_4$$

$$\Rightarrow 2n^3 + 5n^2 + n \text{ FLOPs}$$

$$m_1 := Hx$$

$$m_2 := y - m_1$$

$$m_3 := H^\dagger m_2$$

$$y' := m_3 + x$$

$$\Rightarrow 2n^2 + 2n \text{ FLOPs}$$

Input

n = 1500

m = 1000

Matrix M(n, n) <SPD>

Matrix X(n, m) <FullRank>

ColumnVector y(n) <>

ColumnVector b(m) <>

b = inv(X'*inv(M)*X)*X'*inv(M)*y

Instruction Set

BLAS [DDC⁺90]

- $y \leftarrow Ax + y$
- $C \leftarrow AB + C$
- $B \leftarrow A^{-1}B$
- ...

LAPACK [AB⁺99]

- Matrix factorizations.
- Eigensolvers.
- Solvers for linear systems with specific properties.

Linear Algebra Knowledge

- Properties: trmm vs. gemm

- Inference of properties:  $\begin{matrix} \square & \\ & \square \end{matrix} \begin{matrix} \square & \\ & \square \end{matrix} \rightarrow \begin{matrix} \square & \\ & \square \end{matrix}$

- Simplifications: $A^T \rightarrow A$ if $\text{Symmetric}(A)$

- Rewriting expressions:

$$X := A^T A + A^T B + B^T A \quad \rightarrow \quad \begin{array}{l} Y := B + A/2 \\ X := A^T Y + Y^T A \end{array}$$

- Common subexpressions:

$$X := AB^{-T}C + B^{-1}A^T \quad \rightarrow \quad \begin{array}{l} Z := AB^{-T} \\ X := ZC + Z^T \end{array}$$

- Matrix chains:

$$(AB)c \quad \mathcal{O}(n^3)$$

$$A(Bc) \quad \mathcal{O}(n^2)$$

Derivation Graph

$$y' := H^\dagger y + (I_n - H^\dagger H)x$$

Derivation Graph

$$y' := H^\dagger y + (I_n - H^\dagger H)x$$

$$M_1 := H^\dagger H$$

$$y' := H^\dagger y + (I_n - M_1)x$$

Derivation Graph

$$y' := H^\dagger y + x - H^\dagger H x$$

$$M_1 := H^\dagger H$$

$$y' := H^\dagger y + (I_n - M_1)x$$

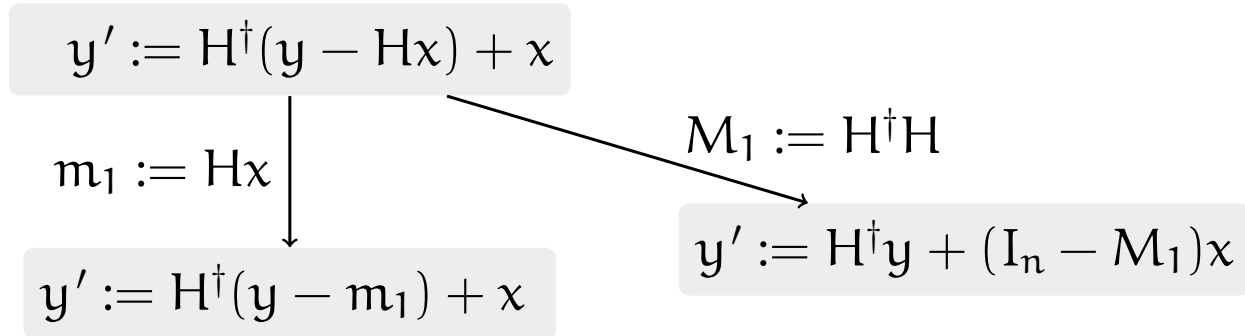
Derivation Graph

$$y' := H^\dagger(y - Hx) + x$$

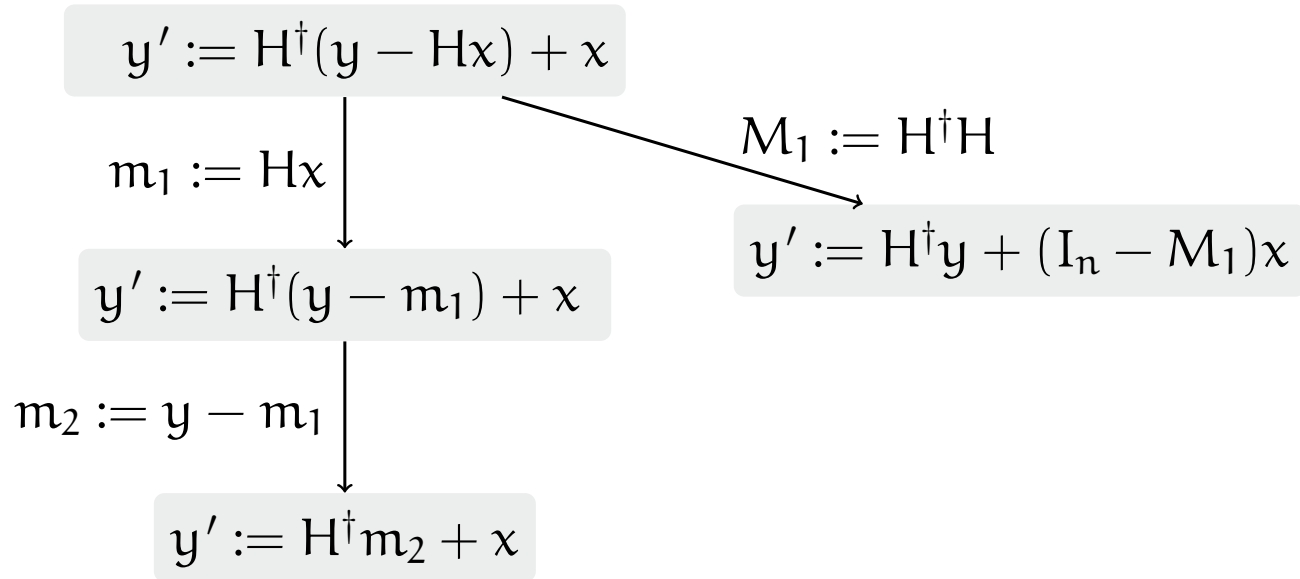
$$M_1 := H^\dagger H$$

$$y' := H^\dagger y + (I_n - M_1)x$$

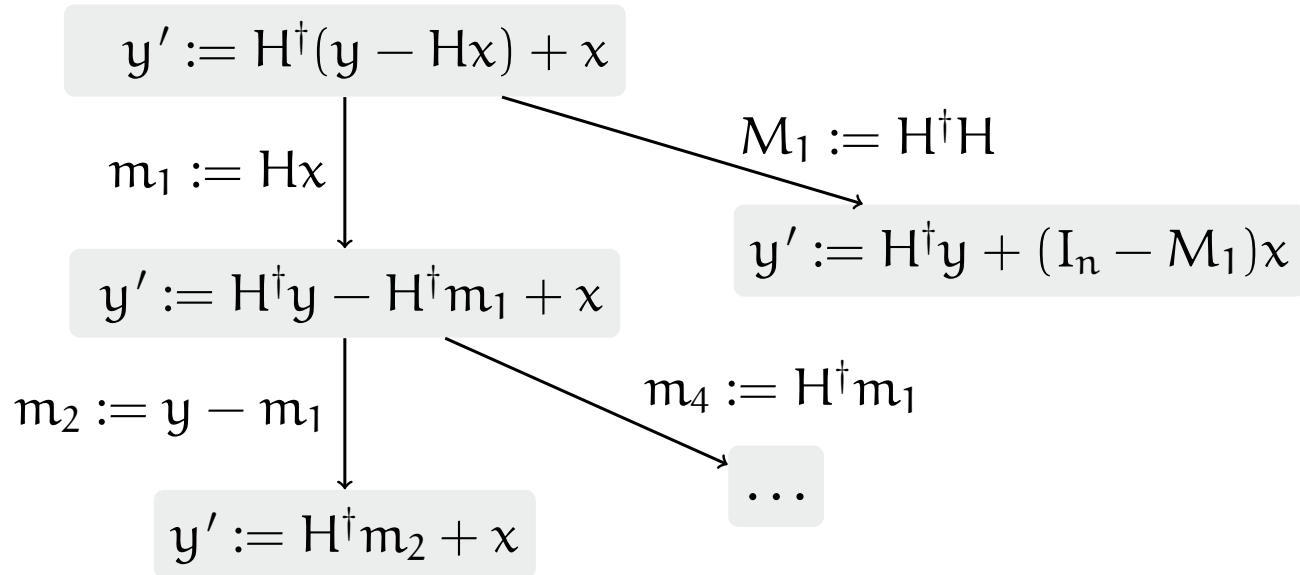
Derivation Graph



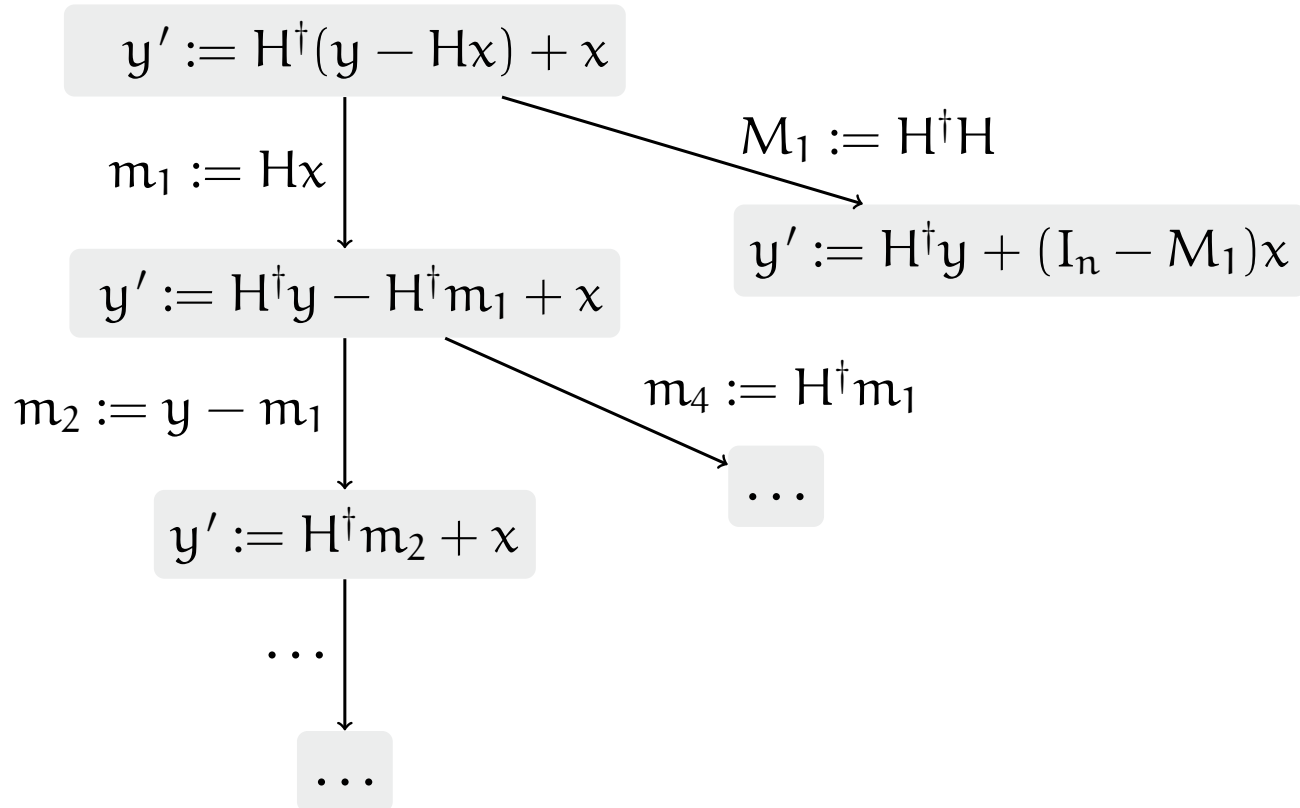
Derivation Graph



Derivation Graph

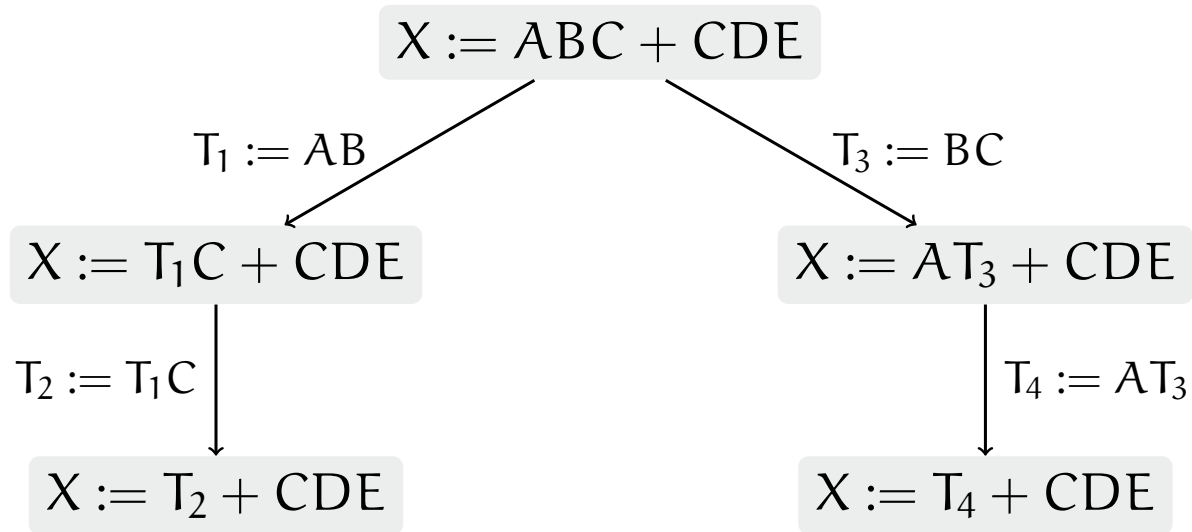


Derivation Graph



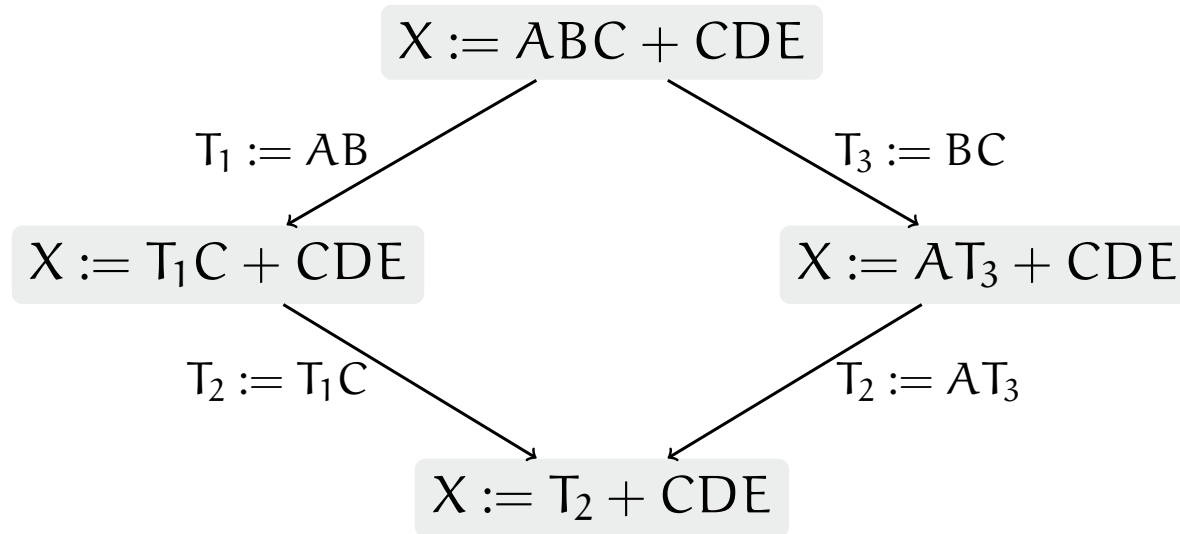
Derivation Graph

Reducing Redundancy



Derivation Graph

Reducing Redundancy



Derivation Graph

Reducing Redundancy

$X := AB + AC + AD$



Derivation Graph

Reducing Redundancy

$X := AB + AC + AD$

$M_1 := AB$

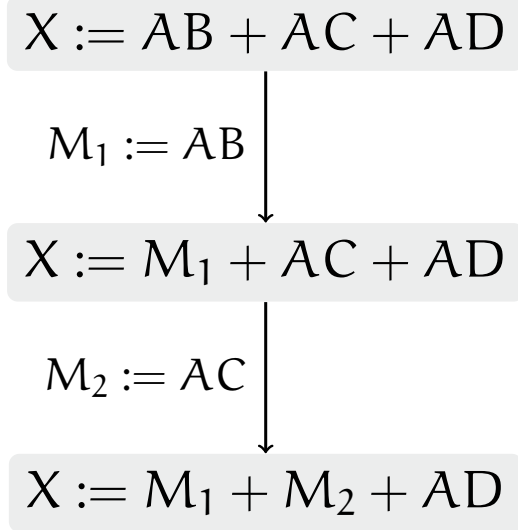


$X := M_1 + AC + AD$

tmp	expr
M_1	AB

Derivation Graph

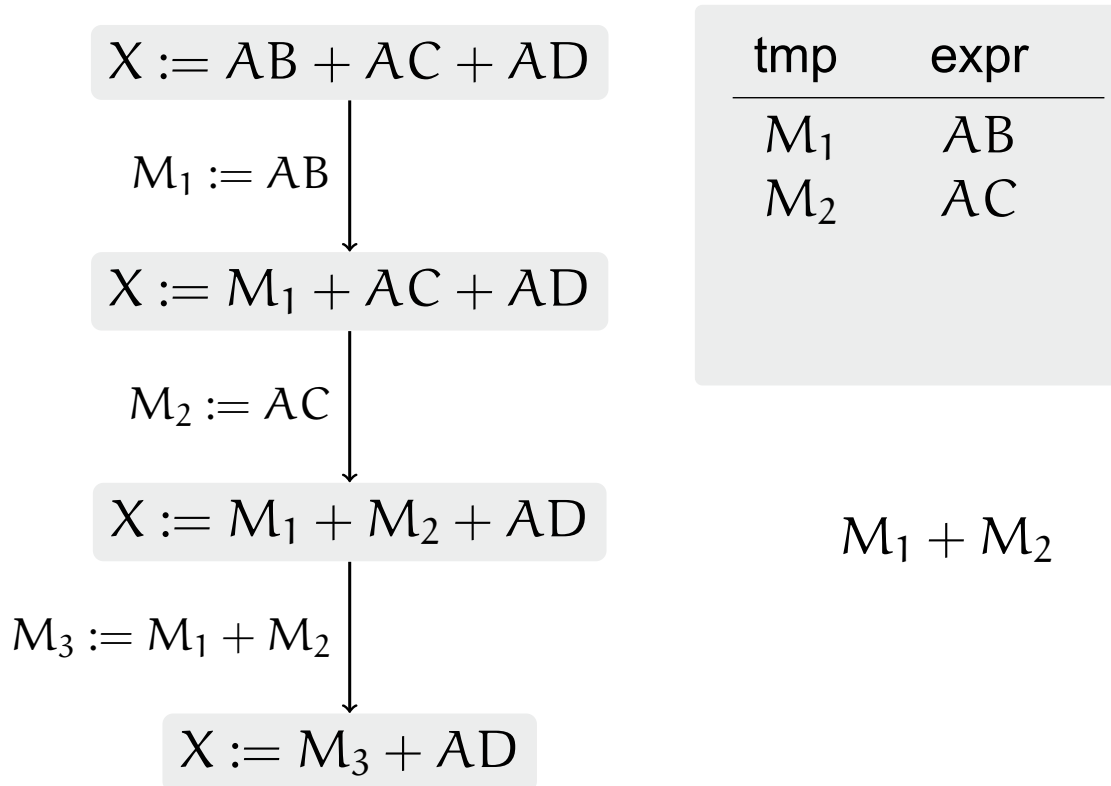
Reducing Redundancy



tmp	expr
M_1	AB
M_2	AC

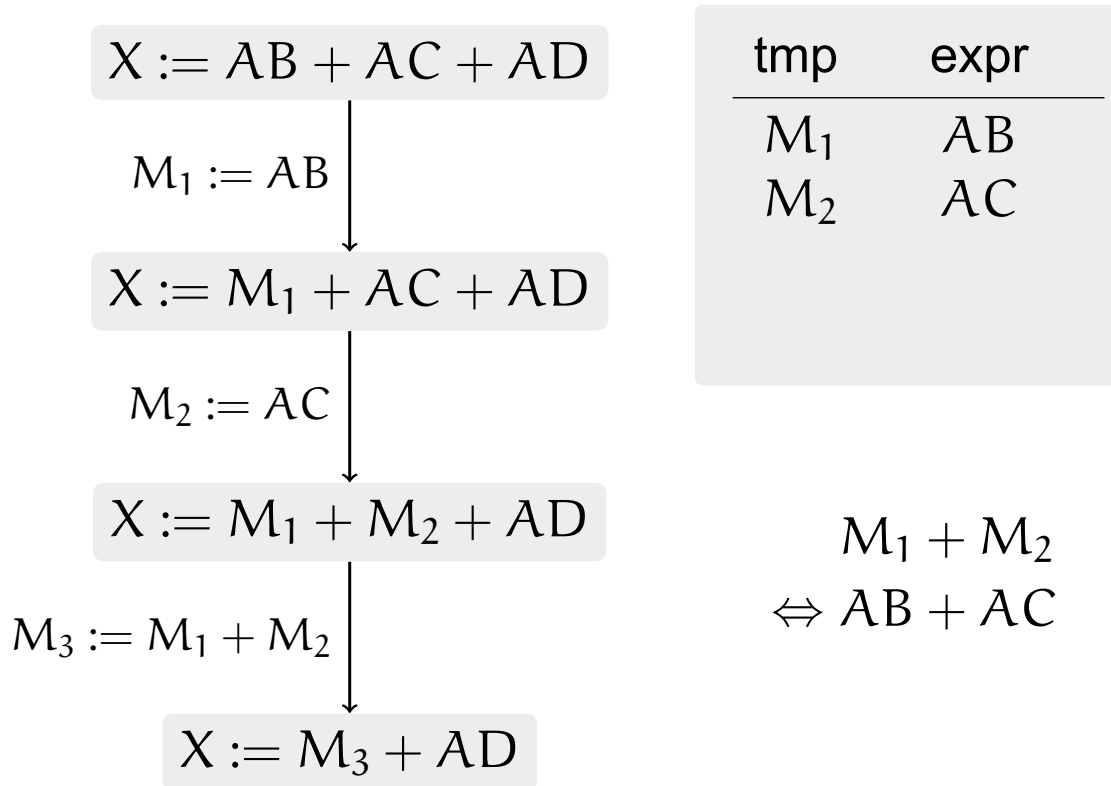
Derivation Graph

Reducing Redundancy



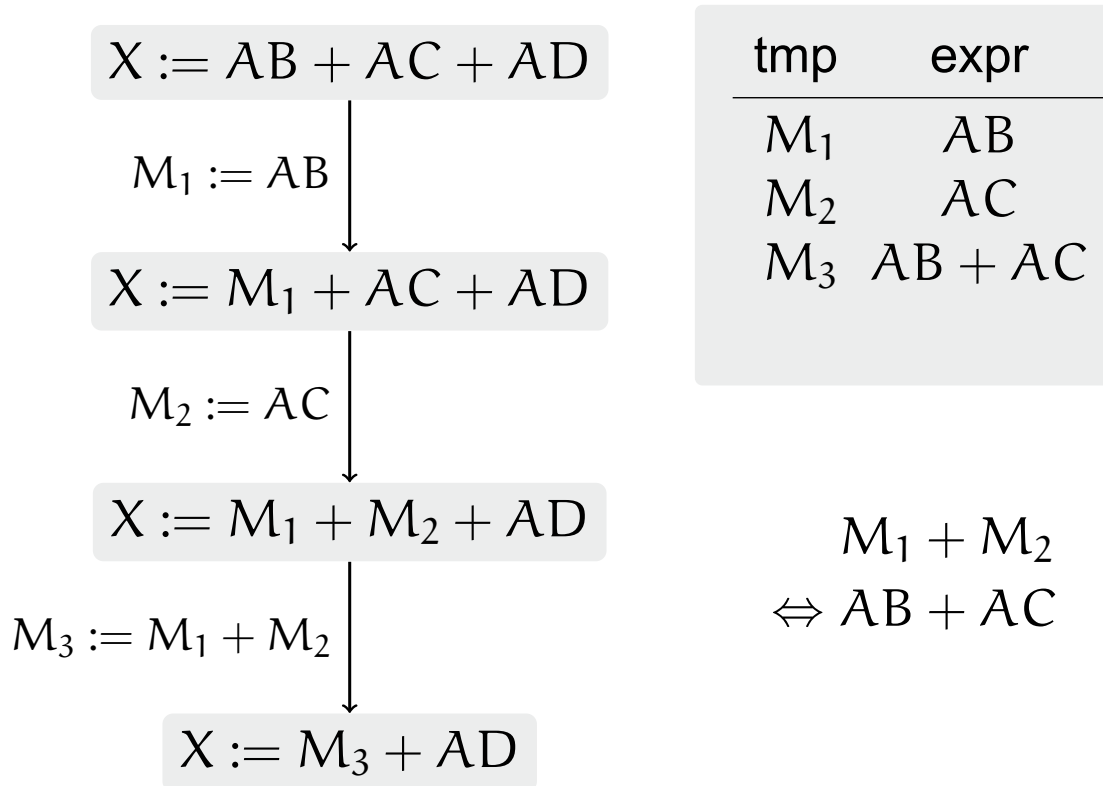
Derivation Graph

Reducing Redundancy



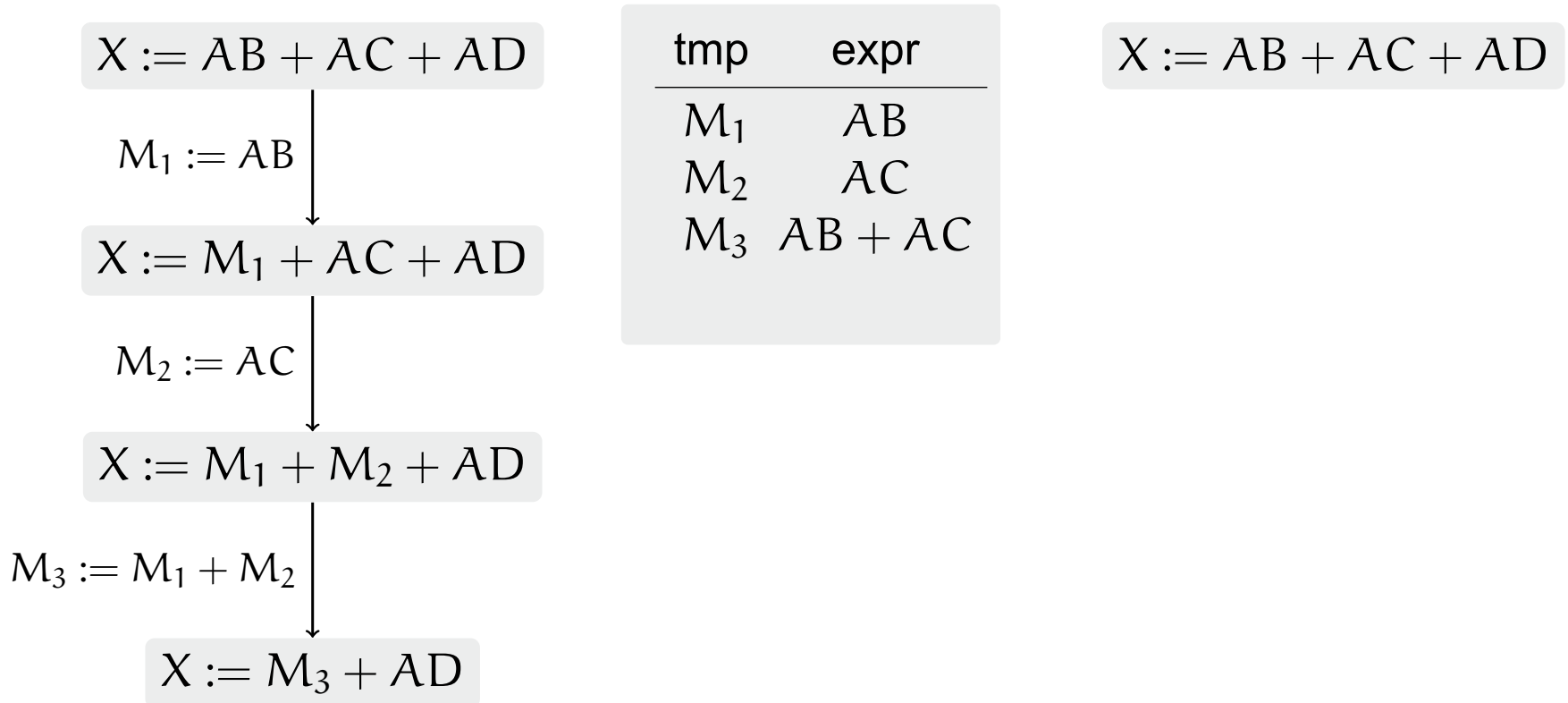
Derivation Graph

Reducing Redundancy



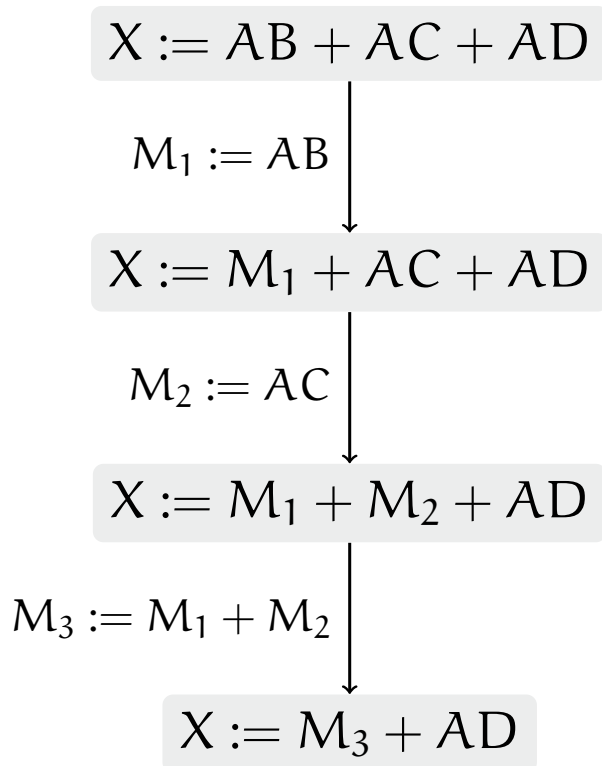
Derivation Graph

Reducing Redundancy



Derivation Graph

Reducing Redundancy

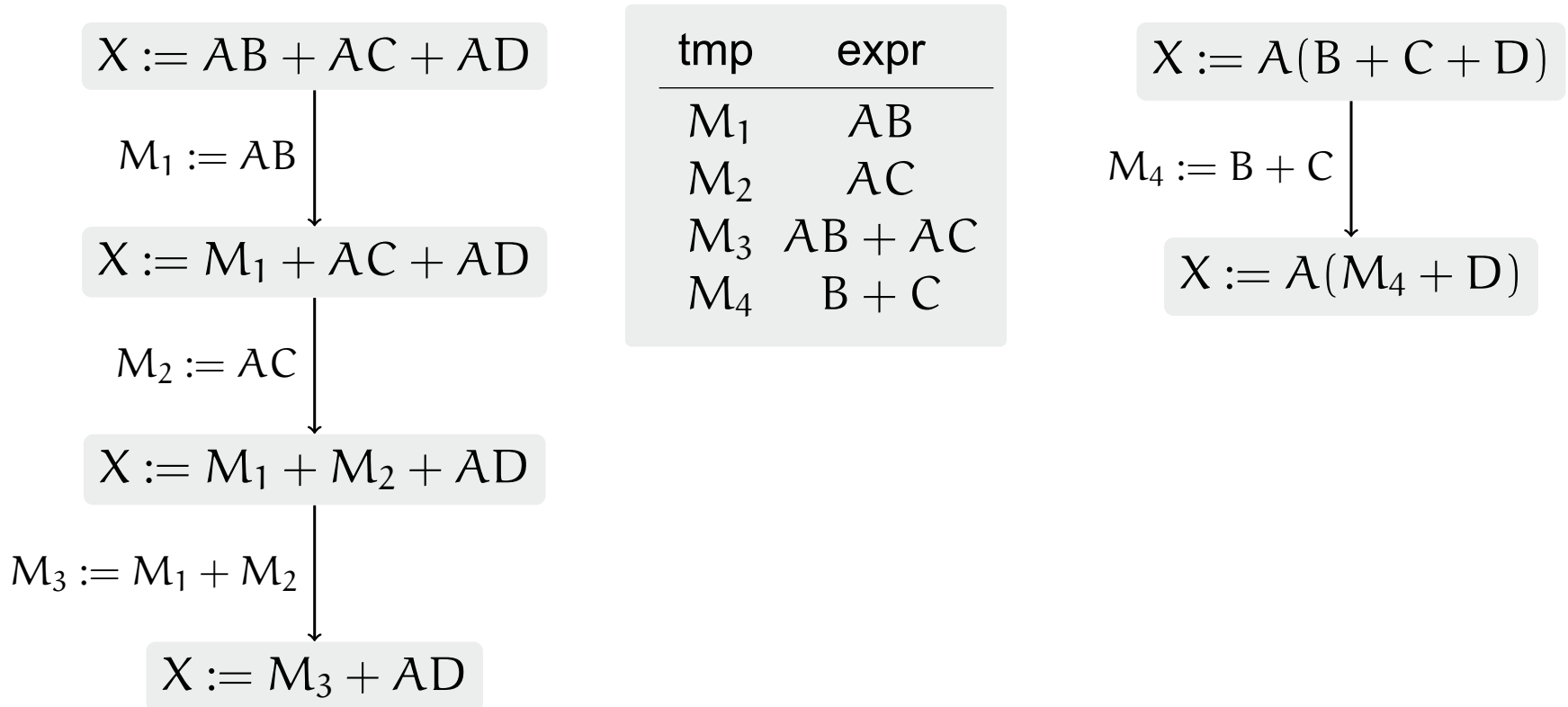


tmp	expr
M_1	AB
M_2	AC
M_3	$AB + AC$

$X := A(B + C + D)$

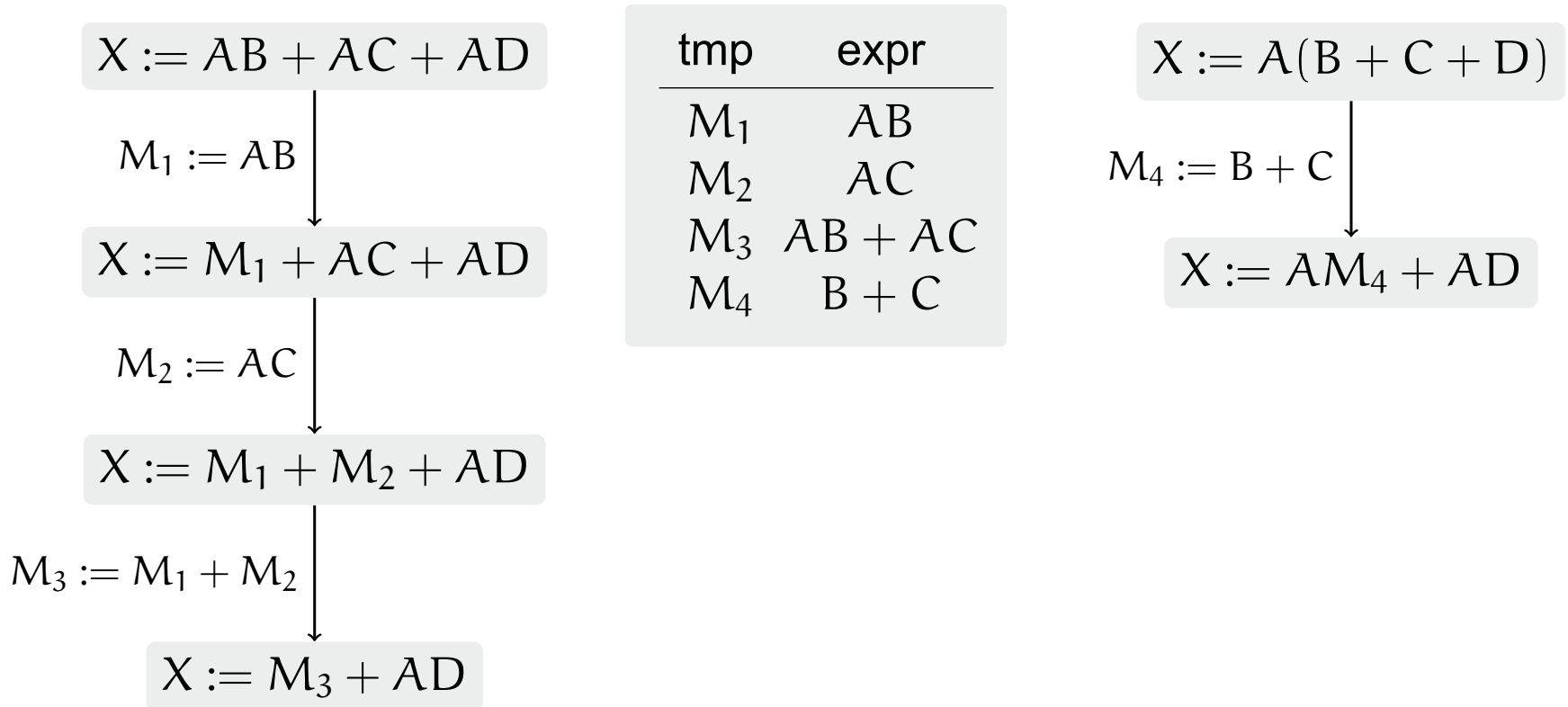
Derivation Graph

Reducing Redundancy



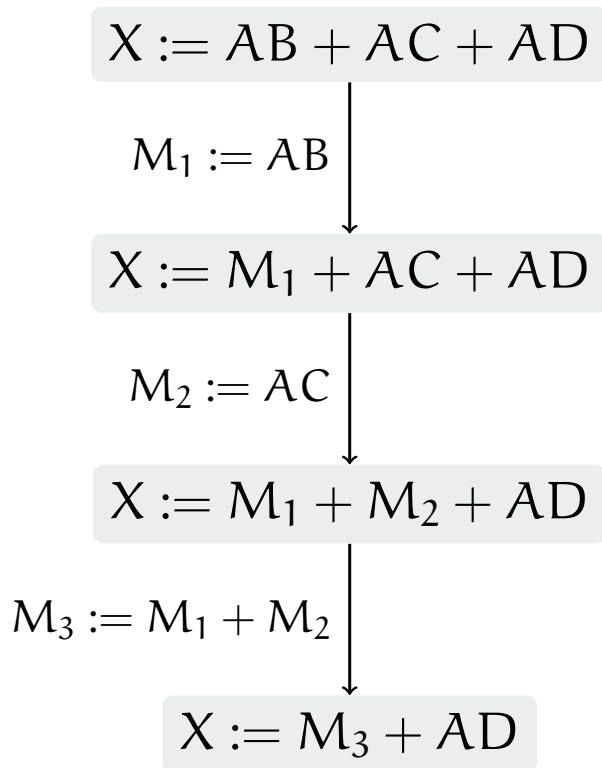
Derivation Graph

Reducing Redundancy



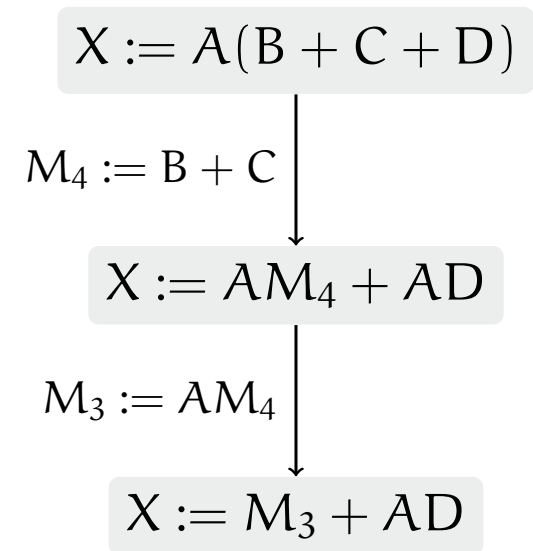
Derivation Graph

Reducing Redundancy



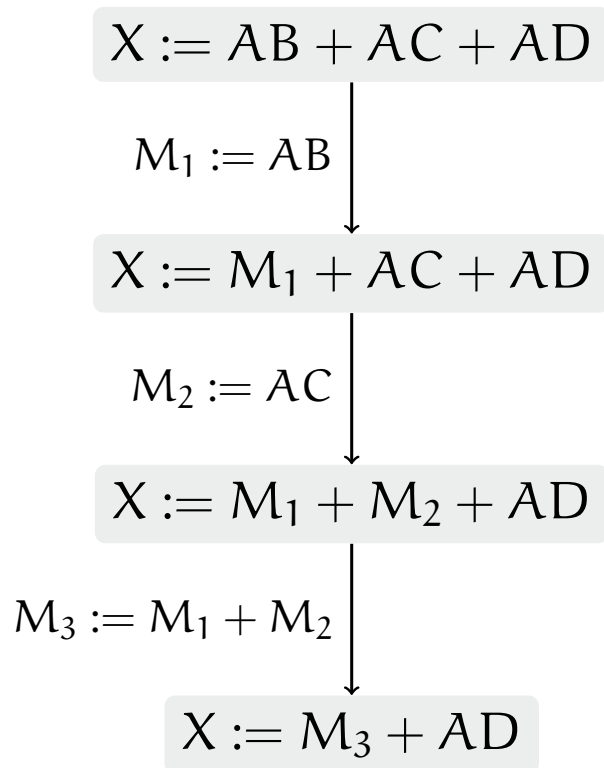
tmp	expr
M_1	AB
M_2	AC
M_3	$AB + AC$
M_4	$B + C$

AM_4



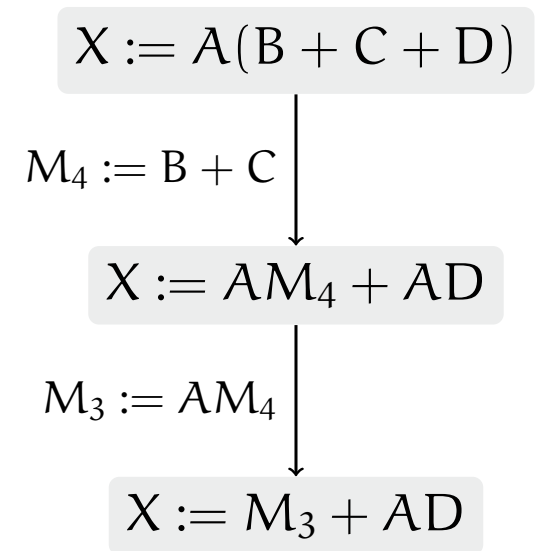
Derivation Graph

Reducing Redundancy



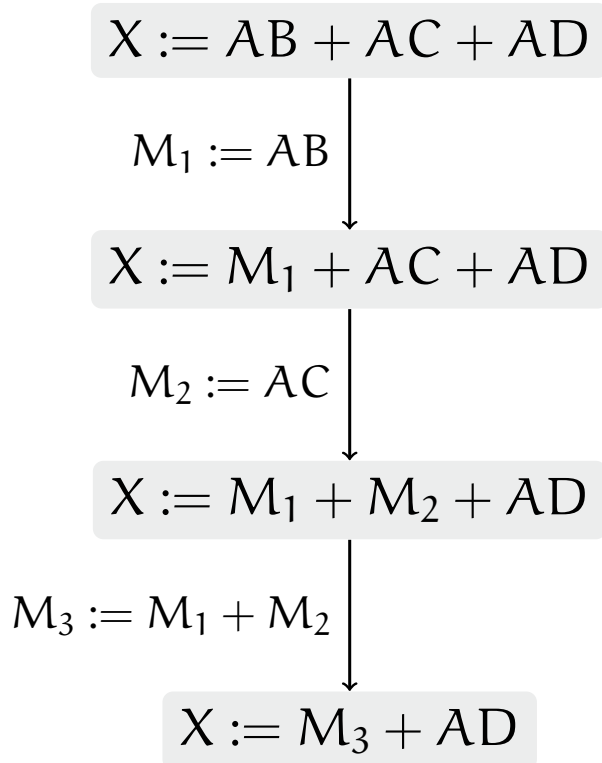
tmp	expr
M_1	AB
M_2	AC
M_3	$AB + AC$
M_4	$B + C$

$$AM_4 \Leftrightarrow A(B + C)$$



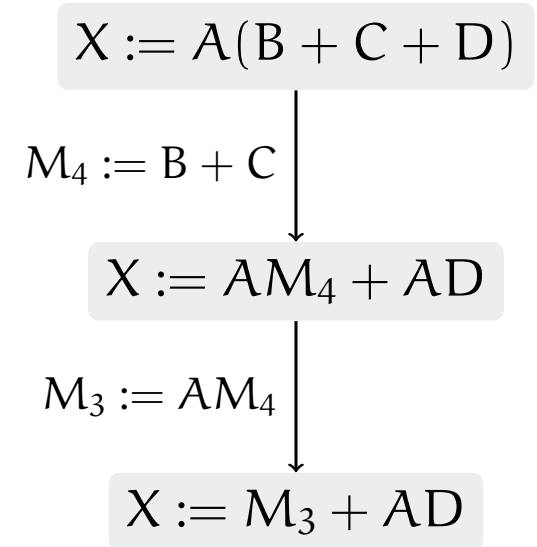
Derivation Graph

Reducing Redundancy



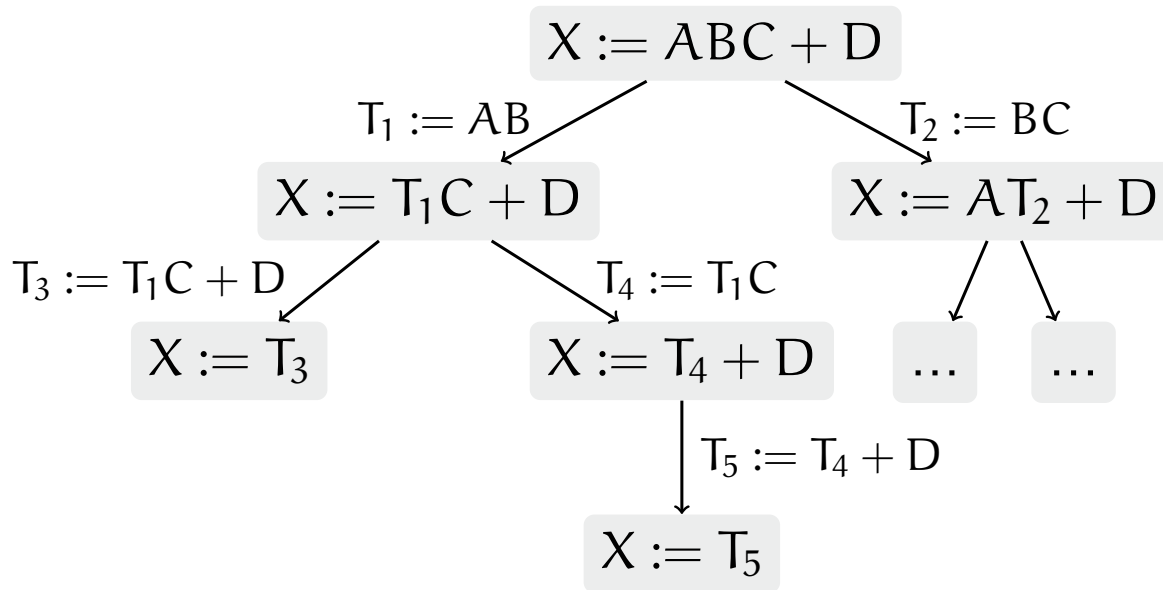
tmp	expr
M_1	AB
M_2	AC
M_3	$AB + AC$
M_4	$B + C$

$$AM_4$$
$$\Leftrightarrow A(B + C)$$
$$\Leftrightarrow AB + AC$$

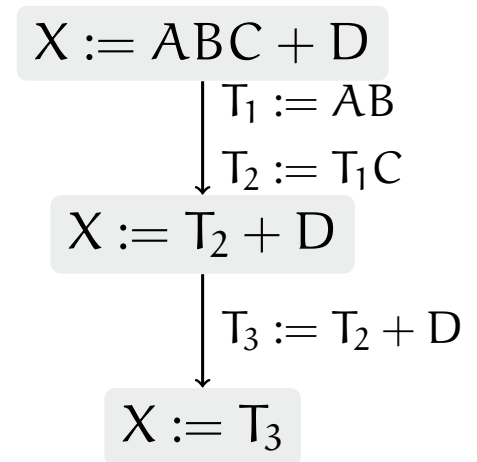


Derivation Graph

Exhaustive



Constructive



Results

Example: $w := AB^{-1}c$

Naive

$w = A * \text{inv}(B) * c$

Recommended

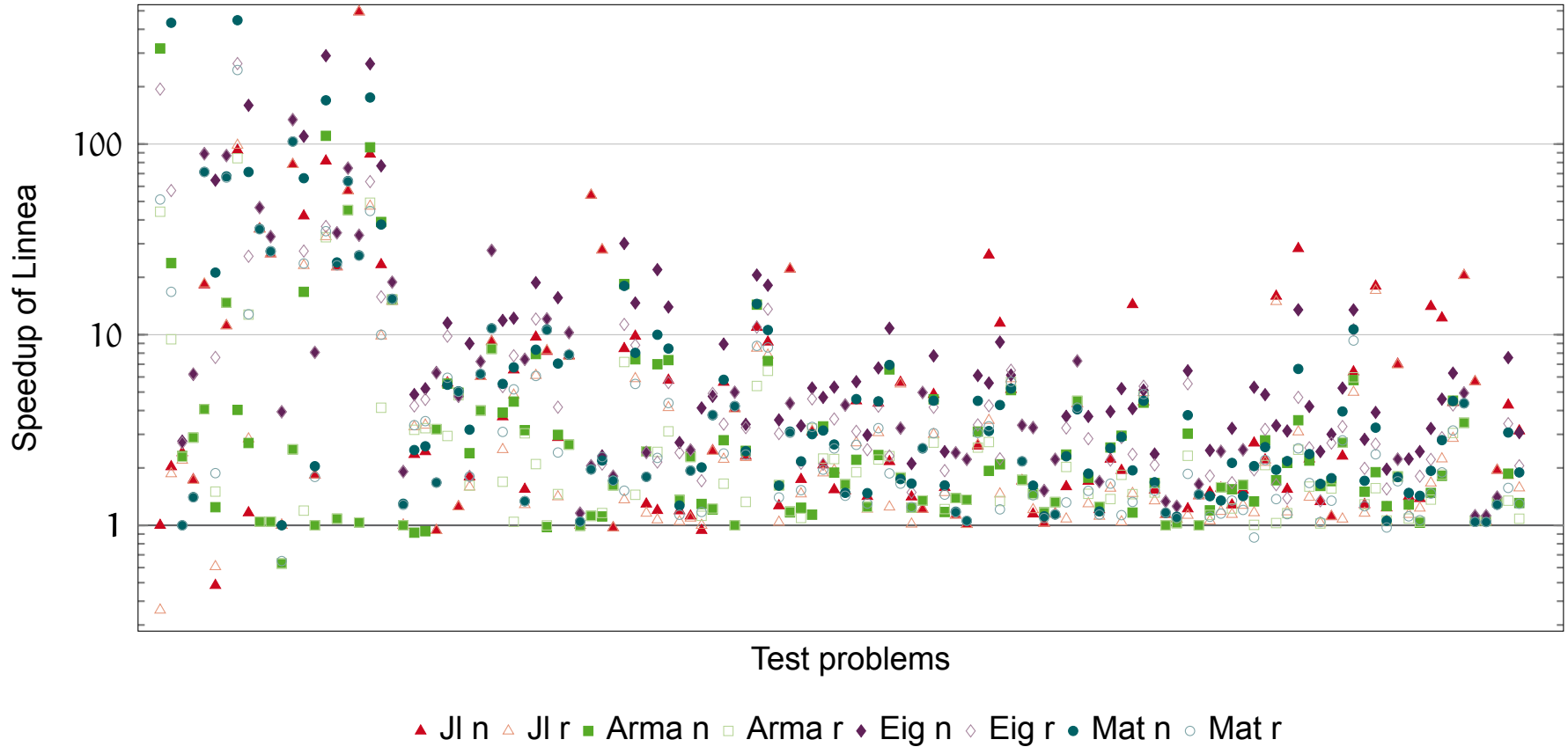
$w = A * (B \setminus c)$

Generated

```
m10 = A; m11 = B; m12 = c;
potrf!('L', m11)
trsv!('L', 'N', 'N', m11, m12)
trsv!('L', 'T', 'N', m11, m12)
m13 = Array{Float64}(1000)
gemv!('N', 1.0, m10, m12, 0.0, m13)
w = m13
```

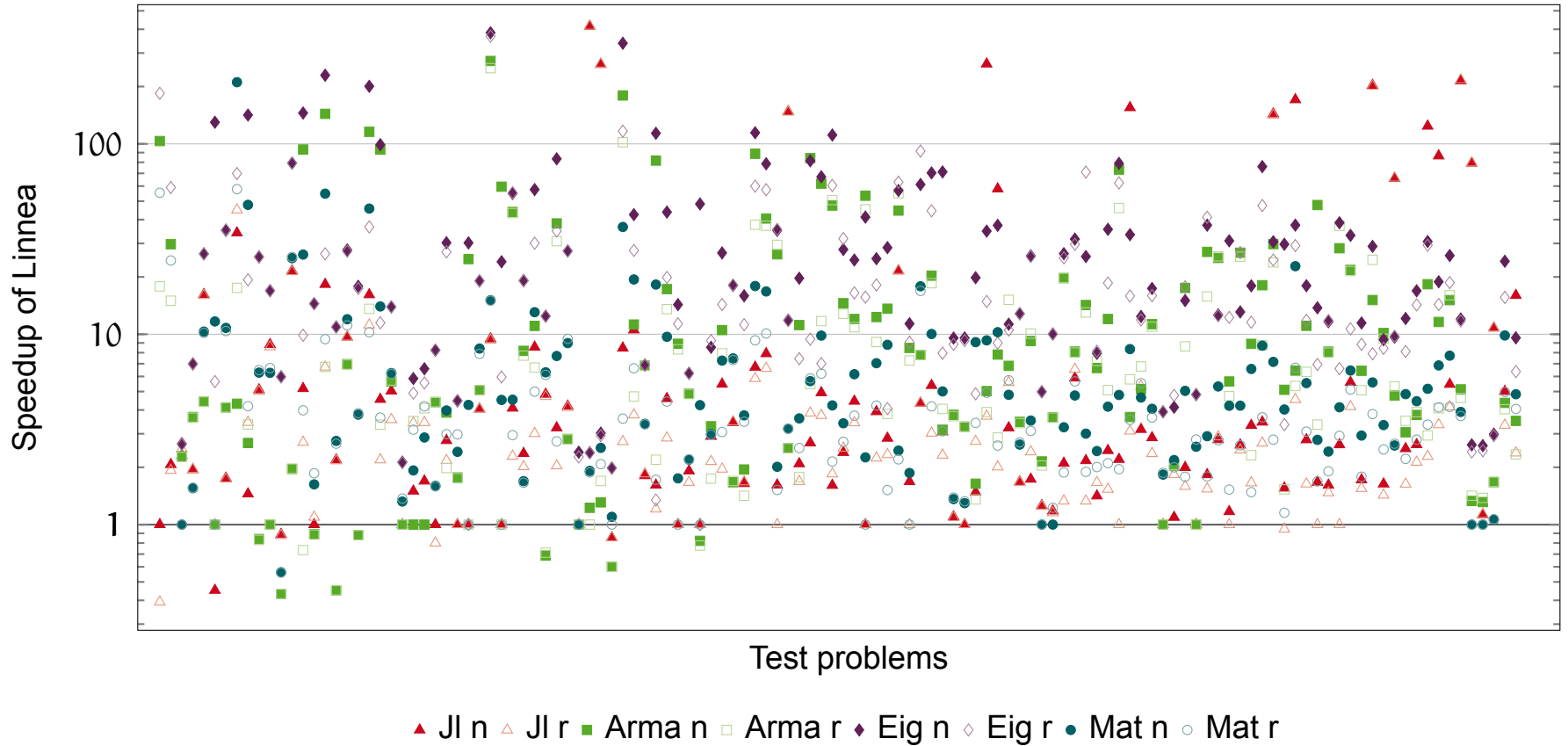
Results

1 Thread



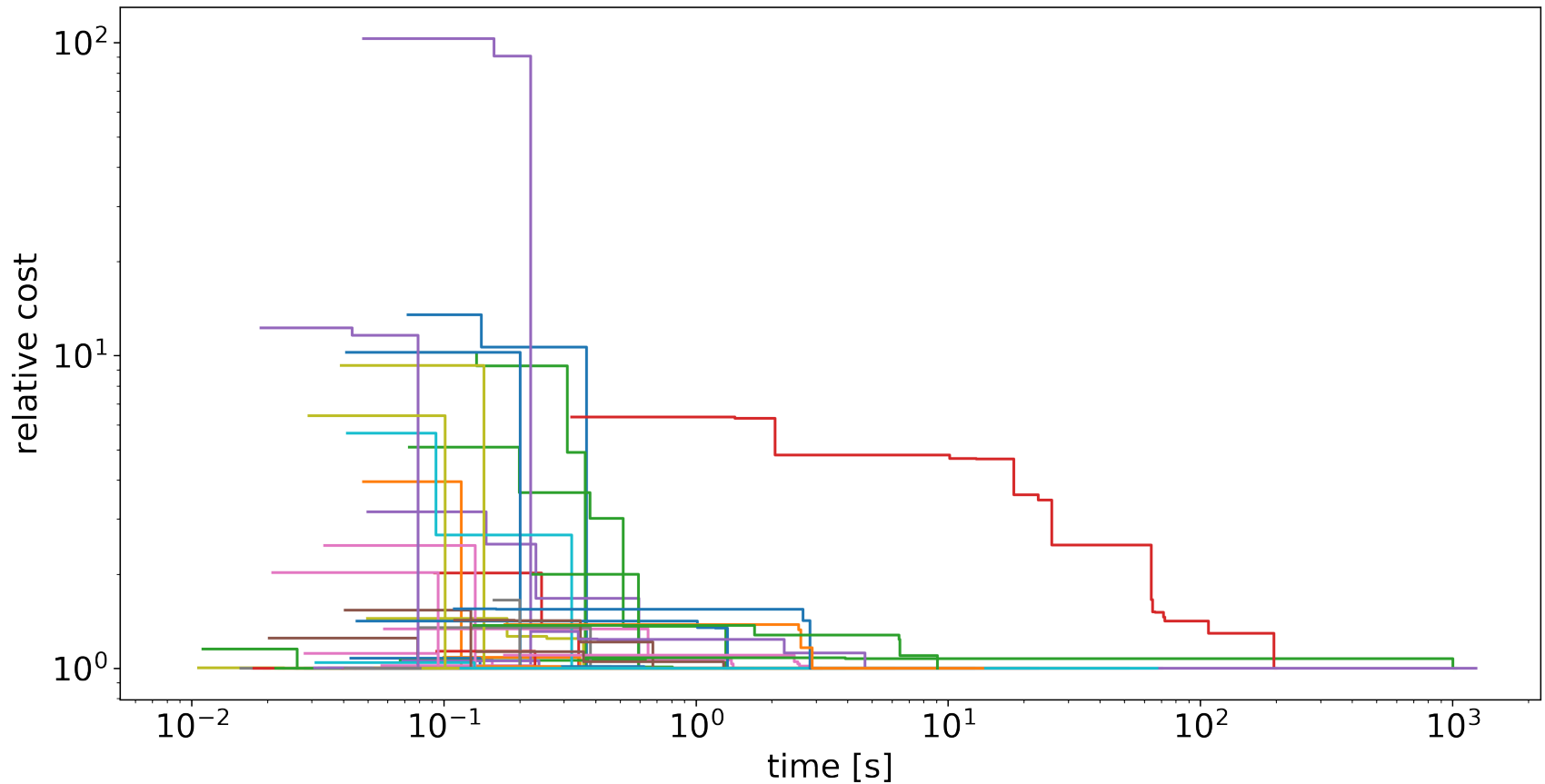
Results

24 Threads



Results

Solutions over Time



Results

- 2x Intel Haswell E5-2680 v3.
 - 24 cores.
 - 64 GB RAM.
 - Turbo Boost is disabled.
- BLAS/LAPACK: Intel MKL 2018.
- Setup:
 - 20 repetitions.
 - We compute the confidence interval of the median [HB15].
 - Cold cache.
- 25 application problems.
 - Domains: statistics, signal processing, image processing, optimization, regularization, linear algebra algorithms.
- 100 randomly generated problems.
 - Between 4 and 7 operands.
 - Sizes between 50 and 2000.
 - Properties: diagonal, lower/upper triangular, symmetric, SPD.
 - No common subexpressions.

References

- [AB⁺99] Edward Anderson, Zhaojun Bai, et al. *LAPACK Users' guide*, volume 9. SIAM, 1999.
- [DDC⁺90] Jack J. Dongarra, Jeremy Du Croz, et al. A set of Level 3 Basic Linear Algebra Subprograms. *ACM TOMS*, 16(1):1–17, 1990.
- [HB15] Torsten Hoefler and Roberto Belli. Scientific Benchmarking of Parallel Computing Systems: Twelve Ways to Tell the Masses When Reporting Performance Results. In *the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 73–12, New York, New York, USA, November 2015. ACM.
- [TG17] Tom Tirer and Raja Giryes. Image Restoration by Iterative Denoising and Backward Projections. *arXiv.org*, pages 138–142, October 2017.

Linnea is available online: <https://github.com/HPAC/linnea>